

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Applicant(s): SHIROGANE, et al.  
Serial No.: Not yet assigned  
Filed: July 30, 2003  
Title: STORAGE SYSTEM  
Group: Not yet assigned

LETTER CLAIMING RIGHT OF PRIORITY

Commissioner for Patents  
P.O. Box 1450  
Alexandria, VA 22313-1450

July 30, 2003

Sir:

Under the provisions of 35 USC 119 and 37 CFR 1.55, the applicant(s) hereby claim(s) the right of priority based on Japanese Patent Application No.(s) 2002-335301, filed November 19, 2002.

A certified copy of said Japanese Application is attached.

Respectfully submitted,

ANTONELLI, TERRY, STOUT & KRAUS, LLP



---

Ronald J. Shore  
Registration No. 28,577

RJS/alb  
Attachment  
(703) 312-6600

日本国特許庁  
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office

出願年月日  
Date of Application:

2002年11月19日

出願番号  
Application Number:

特願2002-335301

[ST.10/C]:

[JP2002-335301]

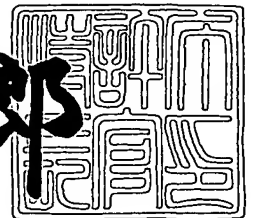
出願人  
Applicant(s):

株式会社日立製作所

2003年 6月 2日

特許庁長官  
Commissioner,  
Japan Patent Office

太田信一郎



出証番号 出証特2003-3041897

【書類名】 特許願

【整理番号】 NT02P0643

【提出日】 平成14年11月19日

【あて先】 特許庁長官 殿

【国際特許分類】 H04B 7/00

【発明者】

【住所又は居所】 神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所 システム開発研究所内

【氏名】 白銀 哲也

【発明者】

【住所又は居所】 神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所 システム開発研究所内

【氏名】 上村 哲也

【発明者】

【住所又は居所】 神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所 システム開発研究所内

【氏名】 田中 淳

【特許出願人】

【識別番号】 000005108

【氏名又は名称】 株式会社日立製作所

【代理人】

【識別番号】 100068504

【弁理士】

【氏名又は名称】 小川 勝男

【電話番号】 03-3661-0071

【選任した代理人】

【識別番号】 100086656

【弁理士】

【氏名又は名称】 田中 恭助

【電話番号】 03-3661-0071

【選任した代理人】

【識別番号】 100094352

【弁理士】

【氏名又は名称】 佐々木 孝

【電話番号】 03-3661-0071

【手数料の表示】

【予納台帳番号】 081423

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 ストレージシステム

【特許請求の範囲】

【請求項 1】

情報発生部からの情報をパケット化して、ネットワークに送信し、またはネットワークからパケットを受信する情報処理装置において、情報発生部からの情報を T C P / I P パケット群に変換する手段と、送信する相手先に対応して F E C (F o r w a r d E r r o r C o r r e c t i o n) の冗長度を管理する管理部と、該管理部に保持している相手先の冗長度を参照して、T C P / I P 変換されたパケット群に対して F E C エンコード処理を行うエンコード部と、ネットワークから受信されるパケット群の情報を F E C デコード処理するデコード部とを有する情報処理装置。

【請求項 2】

前記管理部はメモリに登録されるテーブルを有し、テーブル内に登録される冗長度は相手先ごとに変更可能である請求項 1 記載の情報処理装置。

【請求項 3】

情報処理装置は、情報発生部にディスクドライブを有するストレージ装置であり、かつ情報発生部からの情報を T C P / I P パケットに変換する手段は i S C S I プロトコル処理手段であり、エンコード部は i S C S I のパケット群をエンコードし、デコード部はネットワークから受診したパケット群をデコード処理して i S C S I のパケット群とする請求項 1 の情報処理装置。

【請求項 4】

情報処理装置において、該エンコード部で F E C エンコード処理されたデータは U D P パケット群としてネットワークに送信され、かつネットワークから受信された U D P パケット群は該デコード部で F E C デコード処理される請求項 1 の情報処理装置。

【請求項 5】

ネットワーク側のポートとストレージ装置側のポートとの間に介在し、パケットデータの送受信を中継する装置であって、送信先に対応して F E C 冗長度を登録して管理する送信管理テーブルと、受信先

に対応して F E C 冗長度を登録して管理する受信管理テーブルと、ストレージ装置で生成されるパケット化された iSCSI 層のデータを、該送信管理テーブルを参照して送信先に応じた冗長度を持たせて F E C エンコード処理するエンコード部と、該受信管理テーブルを参照して、ネットワークから受信されるパケットデータを F E C デコード処理し、iSCSI 層のデータに復号化するデコード部とを有する中継装置。

【請求項 6】

前記送信管理テーブルは F E C 処理が可能な送信先アドレスを登録し、かつ前記受信管理テーブルは F E C 処理が可能な送信元アドレスを登録しており、該送信管理テーブルを参照して送信先アドレスが登録されていれば、iSCSI データは前記エンコード部で F E C 処理されてネットワークへ送信され、送信先アドレスが登録されていないければ、iSCSI データは該エンコード部で F E C 処理されずにネットワークへ送信され、一方ネットワークから受信されるパケットデータに対し、該受信管理テーブルを参照して送信元アドレスが登録されていれば該デコード部で F E C デコード処理して iSCSI データに復号化し、登録されていないければ該デコード部で F E C 復号化されずに iSCSI 層へ送る請求項 5 記載の装置。

【請求項 7】

ネットワークを介して受信される制御フレームの内容を解析してアドレスを追加または削除するようにして、前記送信管理テーブルの内容及び前記受信管理テーブルの内容を変更する手段を有する請求項 5 記載の装置。

【請求項 8】

i S C S I プロトコルを用いた装置間でネットワークを介してデータを伝送する通信方法であって、F E C 通信モードでデータの送受信を行う第 1 の通信モードと、T C P / I P 通信モードでデータの送受信を行う第 2 の通信モードと、データの通信先となる相手の iSCSI ネームをメモリに登録して管理するステップと、送信先に対応して F E C の冗長度をメモリに登録して管理するステップと、データを送信する相手先の iSCSI ネームがメモリに登録されているか否かをチェックするステップと、チェックの結果、iSCSI ネームがメモリに登録されていればメモリに登録されている F E C の冗長度に基づいてデータを F E C 処理して該第 1

の通信モードに従って送信するステップと、チェックの結果、iSCSI ネームがメモリに登録されていなければ第2のモードでデータを送信するステップとを有する通信方法。

【請求項9】

請求項8において、更に通信先毎に送信するパケットの喪失率を求めて管理するステップと、パケットの喪失率に従って通信先毎に該メモリに登録されている冗長度を変更するステップとを有する通信方法。

【請求項10】

請求項8において、更に受信側の装置で、送信元に対応してFECの冗長度をメモリに登録して管理するステップと、データを受信するとき、送信元のiSCSI ネームが該メモリに登録されているか否かをチェックするステップと、チェックの結果、iSCSI ネームがメモリに登録されていればメモリに登録されているFECの冗長度に基づいてデータをiSCSIデータに復元処理するステップとを有する通信方法。

【請求項11】

請求項8において、受信側の装置で、iSCSIデータに復元された場合にはACKを送信元に返送し、復元されなかった場合にはACKを送出しないようにするステップと、送信元の装置ではACKを受け取らなかった場合に、先に送信した同じデータを同じ送信先にFEC処理して該第1の通信モードに従って送信するステップとを有する通信方法。

【請求項12】

複数のストレージ装置がネットワークに接続されてデータの送受信を行うストレージシステムにおいて、該ストレージ装置はデータを記録するディスクドライブと、該ディスクドライブに接続されるディスクアダプタと、該ディスクアダプタと接続されるキャッシュメモリと、該キャッシュメモリに接続されるチャネルアダプタと、該ディスクドライブからのデータをTCP/IP上のiSCSIパケット群に変換する手段と、送信する相手先に対応してFEC(Forward Error Correction)の冗長度を管理する管理部と、該管理部に保持している相手先の冗長度を参照して、変換されたTCP/IP上のiSCSIパケット群に対してFE

Cエンコード処理を行うエンコード部と、ネットワークから受信されるパケット群の情報をFECデコード処理するデコード部とを有するストレージシステム。

【請求項13】

データを記録して処理するアプリケーション層と、該アプリケーション層のデータに対してSCSI変換するiSCSI層と、iSCSI層のデータをTCP/IP処理するTCP層及びIP層を有するストレージ装置よりネットワークに対してデータの送受信を行うストレージシステムにおいて、iSCSI層からのデータに対して、冗長符号を付加するFECエンコード処理層と、エンコード処理されたデータにUDP処理するUDP層と、UDP層からのデータに対してIP処理を行うIP層とを有し、該IP層で処理されたデータをネットワークに送出することを特徴とするストレージシステム。

【請求項14】

請求項13において、更に、ネットワークから受信されたデータであって、該IP層でIP処理され、UDP層でUDP処理されたデータをFECデコード処理するFECデコード層を有するストレージシステム。

【請求項15】

請求項13において、該FECエンコード処理層における冗長符号の冗長度をデータの送信先毎に変更する手段を有するストレージシステム。

【請求項16】

iSCSIプロトコルを用いた装置のiSCSI間でネットワークを介してデータを伝送する通信方法であって、FEC処理した通信モードでデータの送受信を行う第1の通信モードと、TCP/IP通信モードでデータの送受信を行う第2の通信モードと、データの通信先となる相手がiSCSI層を有するか否かをiSCSIネームを用いてチェックするステップと、このチェックの結果、相手がiSCSI層を有する場合、相手に応じたFECの冗長度に基づいてデータをFEC処理して該第1の通信モードに従って送信するステップと、チェックの結果、相手がiSCSI層を有しない場合、第2のモードでデータを送信するステップとを有する通信方法。

【発明の詳細な説明】



【 0 0 0 1 】

【発明の属する技術分野】

本発明はストレージシステムに係り、特にネットワークを介して接続される R A I D のようなストレージ装置間、又はストレージ装置とホスト間のデータ通信方法、及びそれを実現させるためのストレージ装置の構成に関するものである。

【 0 0 0 2 】

【従来の技術】

I P ストレージの普及及びデータのリモートコピーに対するニーズの増大により、ネットワークを介して接続された遠距離にあるストレージ間或いはストレージ装置とサーバ間でデータ転送することが頻繁になってきている。しかも、転送するデータ量は年々増加しており、大容量のデータを高速、高信頼で送信することが必要となってきた。

【 0 0 0 3 】

サーバとストレージ装置間を接続するとインタフェースして、ファイバチャネルや S C S I インタフェースが知られている。ファイバチャネルは S A N (Storage Area Network) を構成する高速データ転送のための標準的なインタフェースとして使用されている。また、S C S I インタフェースは、データ伝送速度が速く、伝送遅延も小さい上に、伝送エラーの発生確率が非常に低いと言うメリットがあるので、従来からストレージ装置用のプロトコルとして良く使用されているが、S C S I インタフェースは伝送距離が短いと言う問題がある。

【 0 0 0 4 】

最近、iSCSI プロトコルを用いてストレージ装置を接続する試みが行われている。iSCSI とは、ネットワーク技術である T C P / I P 上で、インタフェース技術として S C S I 処理を実現するためのプロトコル技術であり、現在 I E T F (the Internet Engineering Task Force) で規格化が進められている。この iSCSI プロトコルにより長距離にあるストレージ装置間でデータ転送する場合、I P 層でパケットの喪失が多発することが予想される。

【 0 0 0 5 】

一般にデータ通信におけるパケットの喪失の対策として、A R Q (Automatic R

repeat reQuest)方式即ち自動再送要求方式や、F E C (Forward Error Correction)方式即ち前方誤り訂正方式が知られている。長距離のデータ伝送を行う場合、パケットの喪失が頻繁に生じる可能性がある。そのため、A R Q方式のデータ通信によれば、データ再送のための時間がかかり、結果的にデータの伝送効率が低下する。特に、長距離データ伝送ではデータ再送にかかる時間のうち、ネットワーク遅延が距離に従って増大するため、A R Q方式は好ましいとは言えない。

【 0 0 0 6 】

例えば、特開 2 0 0 1 - 7 7 8 5 号公報(特許文献 1)には、F E C方式とA R Q方式を併用し、伝送路の回線品質に適した誤り制御方式を通信相手に通知するためにユニークワードを付加して送信し、受信側ではユニークワードに応じて複合化する技術が開示されている。また、特開 2 0 0 1 - 1 6 8 9 4 4 号公報(特許文献 2)には、属性の異なる 2 種以上のデータ(即ち I S Oデータと A S Yデータ)を送信し、受信側ではこれらの属性に応じて受信処理する方式が提唱されている。

【 0 0 0 7 】

【特許文献 1】

特開 2 0 0 1 - 7 7 8 5 号公報

【特許文献 2】

特開 2 0 0 1 - 1 6 8 9 4 4 号公報

【 0 0 0 8 】

【発明が解決しようとする課題】

然るに、上記文献 1 及び文献 2 には、ストレージ装置との関係が開示されていないし、また F E C冗長度の変化にどのように対応するのかについても開示がない。

【 0 0 0 9 】

本発明の目的は、ネットワークを介して行うデータ通信においてパケットが喪失した場合でも送信データの復元が可能な通信方法を提供することにある。

本発明の他の目的は、FEC方式を採用した iSCSI プロトコル準拠のストレージ装置を提供することにある。

本発明の他の目的は、iSCSI層間でデータの伝送を行う場合、データの送信先毎にFEC処理の状態やデータ送信の冗長度を変えてデータ送信し、復元することができるストレージシステムを提供することにある。

## 【 0 0 1 0 】

## 【課題を解決するための手段】

上記目的を達成するために、本発明は、ストレージ装置どうし又はストレージ装置とホストコンピュータ（一般的な意味のコンピュータであり、サーバも含む）とがネットワークを介して接続されるストレージシステムにおいて、ストレージ装置はiSCSIプロトコルを採用し、iSCSI層間どうしで夫々データの送受信を行うシステムを達成する。ストレージ装置はFEC制御ユニットを有し、送信側ではこのFECユニットでデータを冗長化（エンコード処理）し、この冗長化したパケット群をネットワークに送信する。一方受信側ではネットワークから受信した冗長パケット群を元にFECユニットで受信データを復元（デコード処理）する。このようにFEC冗長化処理及び復号化処理を行うことにより、冗長化したパケット群の一部が喪失して受信できなくても、元のパケット群は復元可能である。一方受信側でデータの復元できなかった場合には、送信側にACKが返送されないので、送信側では例えばタイムアウト監視を行うことによりデータを再送することができる。

## 【 0 0 1 1 】

本発明は送信するパケットに対するFECの冗長度を変化させることができる。通信する相手のパケット喪失の状態に従って送信側でのFECの冗長度の状態を変える。冗長度の変更をする手法としては、例えば送受信するデータに対して冗長化を行ったデータ量の比率を変える等のやり方がある。また、パケットの喪失率は、例えば受信側で受信するパケット群の喪失数を送信相手毎に採取して求めることができる。これに基づいて送信側ではデータの冗長度を変えて送信する。例えば、パケットの喪失率が高ければ、FECの冗長度を高くして送信し、一方喪失率が低ければ、冗長度を下げて送信する。ネットワークによってパケットの喪失率が時々変化することが予想されるが、それに応じて冗長度を変えられる。また、パケットの喪失率に従って、パケット送信の間隔を変えることも可能であり、広

義にはパケット喪失率に従ってデータ送信の状態を変えるようにできる。

本発明は、好ましい例においては送信先に対応してFEC冗長度を登録して管理する送信管理テーブルと、受信先に対応してFEC冗長度を登録して管理する受信管理テーブルと、ストレージ装置で生成されるパケット化されたiSCSI層のデータを、送信管理テーブルを参照して送信先に応じた冗長度を持たせてFECエンコード処理するエンコード部と、受信管理テーブルを参照して、ネットワークから受信されるパケットデータをFECデコード処理し、iSCSI層のデータに復号化するデコード部とを有して構成される。これは好ましい例では、ストレージ装置に接続されるアダプタとして構成されるが、ストレージ装置内に組み込まれて構成されるようにしても良い。

#### 【0012】

本発明はまた、iSCSIプロトコルを用いたストレージシステムのiSCSI I 間でネットワークを介してデータを伝送する通信方法を実現できる。FEC処理した通信モードでデータの送受信を行う第1の通信モードと、TCP/IP通信モードでデータの送受信を行う第2の通信モードと、データの通信先となる相手がiSCSI層を有するか否かをiSCSI ネームを用いてチェックするステップと、このチェックの結果、相手がiSCSI層を有する場合、相手に応じたFECの冗長度に基づいてデータをFEC処理して該第1の通信モードに従って送信するステップと、チェックの結果、相手がiSCSI層を有しない場合、第2のモードでデータを送信するステップとを有する。

iSCSI層間で通信するストレージシステムの好ましい例によれば、TCPセッションを開始するiSCSIのログインフレームの送信を監視し、iSCSI ネームを取り出し、それが予め決められた通信相手(iSCSIノード、例えばiSCSIイニシエータやiSCSIターゲット)へのセッションである場合、このセッションが存在する間、その目的アドレスへ送信するときは冗長度を上げるように変化させて送信する。

#### 【0013】

##### 【発明の実施の形態】

以下、図面を参照して本発明の実施例を説明する。

図1は本発明の一実施例によるネットワークに接続されたストレージシステムを

示すブロック図である。図 1 において、I P ネットワーク 4 0 0 にはストレージ装置 100, 200、及びホストコンピュータ(単にホストと言う)300が接続されている。このシステムにおいて、ネットワーク 400 を介してストレージ装置 100, 200 間、又はストレージ装置 100, 200 とホスト 300 間でパケット形式のデータの送受信が行われる。ストレージ装置間の通信は例えばリモートコピーを行う場合であり、ストレージ装置とホスト間の通信は例えばデータ処理時、或いはデータステーションとして使用される場合である。特徴的なことはストレージ装置 100, 200 及びホスト 300 は夫々ネットワークとの接続口に F E C 変換アダプタ 110, 210, 310 を具備することである。

## 【 0 0 1 4 】

各ストレージ装置 100, 200 は、通常多数のディスクドライブ 101, 201 と、ディスクドライブ 101, 201 に S C S I インタフェース 102, 202 を介して接続されるディスクアダプタ 103, 203 と、これらのディスクアダプタ 103, 203 にバス 104, 204 を介して接続されるキャッシュメモリ 105, 205 と、キャッシュメモリ 105, 205 にバス 106, 206 を介して接続されるチャネルアダプタ 107, 207 を具備して構成される。チャネルアダプタ 107, 207 は夫々ポート 108, 208 からギガビットイーサネット(登録商標)のような高速 I P インタフェース 109, 209 を介して F E C 変換アダプタ(Forward Error Correction)110, 210 に接続され、これら F E C 変換アダプタ 110, 210 は高速インタフェース 111, 211 を介して I P ネットワーク 400 に接続されている。チャネルアダプタ 1 0 7, 2 0 7 は夫々 i S C S I プロトコルの処理を行う。尚、図示していないが、高速 I P インタフェース 111, 211 はポートと介して I P ネットワーク 400 に接続される。

## 【 0 0 1 5 】

また、ホスト 300 は、情報発生部としての情報処理装置 301 は演算処理装置又はメモリを有し、これらは内部バス 302 を介してホストバスアダプタ(HBA)303 に接続されている。HBA 303 はポート 304 から高速 I P インタフェース 305 を介して F E C 変換アダプタ 310 に接続されている。F E C 変換アダプタ 310 は I P インタフェース 311 を介して I P ネットワークに接続される。ホスト 300 では、iSCSI HBA 303 がホストの入出力処理を行い、また、ホストの O S (Operating System) から

指示を受けてコマンドを発行する。

上述のストレージ装置100、200およびホスト300に接続されるFEC変換アダプタ110,210,310でFEC符号のエンコード、及びデコードが行われる。このアダプタの詳細な構成は後述する。

尚、図1には示していないが、IPネットワーク400には本発明の対象となるFEC変換アダプタ110,210,310を備えていない従来のストレージ装置やホストが接続されていることがある。

#### 【0016】

図2は本発明の実施例におけるストレージ装置に接続されるFEC変換アダプタの内部構成を示す図である。

例えばFEC変換アダプタ110は、送信系と受信系の双方の系列から構成される。即ち送信系は、ストレージ装置100側の物理層1101、ストレージ装置からのパケットデータを一時格納するバッファ1102、送信パケットデータをFECエンコード処理するFECエンコード部113、FEC処理されたデータを一時格納する送信バッファ1103、及び物理層1104から構成され、物理層1104は高速IPインタフェースを介してネットワークに接続される。

#### 【0017】

一方、受信系は、物理層1114、受信データを一時格納するバッファ1113、FECデコード部115、バッファ1112、及び物理層1111から構成される。FECデコード部115では、受信されたパケットデータをデコード処理してFEC復号化する。

#### 【0018】

またFEC変換アダプタ110は、制御プロセッサ116と、FEC送信管理テーブル117、及びFEC受信管理テーブル118と、FEC許可テーブル119を有している。制御プロセッサ117はアダプタ110内の全体の制御を行う。送信管理テーブル117はこのストレージ装置100から送信可能な相手先のストレージ装置またはホスト等のアドレスが登録されている。FEC許可テーブル119はFEC通信を許容するか否かを管理するテーブルであり、通信相手を特定するコードを登録する。尚、これらテーブルの構成は図7を参照して後述される。

## 【 0 0 1 9 】

送信管理テーブル117に送信先アドレスが登録されていれば、送信データはF E Cエンコード部113でF E C処理されてネットワークへ送信され、一方送信先アドレスが登録されていなければ、データはF E Cエンコード部113でF E C処理されずにネットワークへ送信される。送信データとなるiSCSIデータのF E C処理は、このテーブルに登録されている送信先毎の冗長度に応じて行われる。また冗長符号としては、例えば公知のX O R演算によるパリティビット符号又はReed Solomon符号等を用いることができる。

## 【 0 0 2 0 】

受信管理テーブル118は送信元となるストレージ装置等のアドレスが登録されている。この受信管理テーブル118に送信元アドレスが登録されていればF E Cデコード部115でF E Cデコード処理して内部データであるiSCSIデータに復号化し、登録されていなければデコード部115でF E C復号化されずにiSCSI層へ送られる。

F E C許可テーブル119は通信相手を制御するための台帳としてのテーブルである。これには通信相手のiSCSIネーム毎に冗長度が登録される。

## 【 0 0 2 1 】

図3は本発明の他の実施例を示すストレージシステムのブロック図である。この例では、ストレージ装置100' ,及び200' 共に、ディスクドライブ101,201、ディスクアダプタ103,203、キャッシュメモリ105,205、チャネルアダプタ107,207に至るまで信頼性を考慮して全て二重化されている。これに合せて、F E C変換アダプタ110,210も夫々二重化されて、ネットワーク400に接続されている。これら各々のF E C変換アダプタ110,210がエンコード機能とデコード機能を具備している。この例では、高速インタフェース220にはF E C変換アダプタが接続されていない従来のポートと考えてよい。

## 【 0 0 2 2 】

更に他の例として、ネットワーク500には管理サーバ500が接続されても良い。管理サーバ500は、各F E C変換アダプタ110,210,310が持つF E C送信管理テーブル117, F E C受信管理テーブル118,及び許可テーブル119の管理を行うもので

あり、各通信先（例えばポート）の追加、削除を指示することにより、各 F E C 変換アダプタ間の通信に F E C を使用するか否かを制御する。管理テーブルは F E C の通信を行う F E C 変換アダプタのアドレス、またはポートのアドレスを登録することによりその管理を行う。勿論、F E C 変換アダプタを備えない従来型のストレージ装置やホストは対象とするアドレスがこのテーブルに登録されないで、F E C の効果を享有できない。

## 【 0 0 2 3 】

図 4 は本発明の一実施例によるデータ転送の概念を示す図である。ストレージ装置 100 とストレージ装置 200 がルータ A、B 及びネットワークを介して遠距離でデータを転送する場合を示している。ここで、ストレージ装置 200 はそれに限らず、F E C 変換アダプタを備えるホストまたはサーバであっても良い。

## 【 0 0 2 4 】

ストレージ装置 100 のアプリケーション層で処理されたデータは iSCSI 層でプロトコル変換されて、更に T C P 層、I P 層での制御情報を付加してパケット処理されて、ネットワークインタフェース (I F) を介して F E C 変換アダプタ 110 に送られる。

F E C 変換アダプタ 110 では、送信データに対してエラー訂正のための冗長符号を付加する F E C 処理（エンコード処理）する。このエンコード処理については後で詳述される。エンコード処理の後 U D P ヘッダのポート番号を付加し、更に I P パケット処理して、ネットワーク I F を介してルータ A に送る。データはルータ A の I P 層、ネットワーク I F を介してネットワークに送出される。

## 【 0 0 2 5 】

一方、受信側でデータはルータ B のネットワーク I F、I P 層を介して受信され、F E C 変換アダプタ 210 へ送られる。受信データはネットワーク I F、I P 層、U D P 層を介して F E C 層に送られ、そこで F E C の復号化処理（デコード処理）が行われる。復号化処理においては後で詳しく説明されるが、F E C の冗長符号で訂正処理される。エラーの度合いが大きい場合には A C K が中々送信側に返送されない。そこで、送信側では A C K 受領のタイムアウトを監視することにより、もし既定時間内に A C K を受信しなければ送信側から同じデータを再



送する。

復号化されたデータはストレージ装置200に送られ、同様に各層を介してiSCSI層、アプリケーション層へと送られ、使用に供される。

#### 【 0 0 2 6 】

図5は通信に用いられるパケットのフォーマットの一例を示す図である。

iSCSI層においては、データの通信の単位となるiSCSI PDU(Protocol Data Unit)は、B H S (Basic Header Segment)とデータセグメントから成る。尚、B H S とデータセグメントとの間にA H S (Additional Header Sequence)が入る場合もあるが、図示の例では省略してある。B H S はメッセージの長さを格納しており、データセグメントの開始位置やメッセージの境界が分かる。iSCSI層において、イニシエータとターゲットはiSCSI PDUと言うメッセージで通信を行う。iSCSI PDUの長さは4バイトの倍数である。

T C P 層、I P 層、データリンク層では上記のiSCSI層からのパケットデータに対して先頭に、データリンクヘッダ(D L H)、I Pヘッダ(I P H)、T C Pヘッダ(T C P H)が付加される。一方iSCSIパケットデータの最後部にはデータリンクトレイラ(D L T)が付加される。尚、イーサネットの場合にデータリンクヘッダはイーサヘッダになる。

#### 【 0 0 2 7 】

図2のバッファ(エンコードバッファ)1102にはiSCSIのパケットデータに各ヘッダが付加されたものが順次格納される。そして、I Pヘッダ、T C Pヘッダ等も含めて冗長化される。冗長化されたデータには、データリンクヘッダ(D L H)、I Pヘッダ(I P H)、T C Pヘッダ(T C P H)の他に、更にU D Pヘッダ、F E Cヘッダが付加されて、送信される。

#### 【 0 0 2 8 】

本発明においては、F E C通信を行う場合I Pアドレスではなく、WWN(world wide name)即ちiSCSIネームをF E C許可テーブルに登録し、それによって通信先の指示等を行う。iSCSIログインを監視してF E C許可テーブルに登録されている通信先とのログインが行われたら、前述したF E C送信管理テーブル及びF E C受信管理テーブルに通信先のアドレスを登録することによりF E C通信を

開始する。IPアドレスは例えばポート毎に付与されるものであるが、iSCSIネームは物理的な1つのストレージ装置に複数のネームが設定されることがある。例えば、1つのストレージを複数のパーティションに分割して複数のディスクとして使用する場合には、1つの物理的なポートに対して複数のiSCSIネームが設定される。この利用の仕方はストレージ装置の利用者にとっては大変有意義である。

このように、通信先について例えばIPアドレス毎ではなく、iSCSIネーム毎に管理を行う理由は、iSCSIネームの方が例えば、

“iqn.1993-11.com.disk-vender.diskarrys.sn.45678”のように、ストレージ装置の利用者にとって覚え易く、意味を持たせ易いこと、及びIPアドレスを管理する場合にはネットワークIFカード（又はネットワークアダプタカード）を交換したり、接続先のネットワークの設定が変更になったとき、IPアドレスを設定し直さねばならないが、iSCSIネームではその必要が無いこと、及びIPアドレスはiSCSIネームに比べ詐称される可能性が高いこと等が上げられる。

#### 【 0 0 2 9 】

図6はFECヘッダのフォーマットを示す図である。

FECヘッダは32ビット、4ワードの構成であり、ワード0はエンコード情報、FEC情報の種類、及びFEC制御情報IDからなる。制御情報の種類により異なるワード1、2、3のFEC制御情報の内容が変る。ここで、エンコード情報は、FECデータ部がFEC冗長化されているか否かを示す。FEC情報の種類は、FEC許可テーブルの情報の変更、制御テーブル情報の変更、FEC\_ACK、パケット喪失率の報告、データ、冗長データの種類を示す。FEC制御情報IDはFEC\_ACK、FEC\_RJTが対応する制御情報を示すために用いる。FEC制御情報には、データ長、冗長度、パケット到達率、等の情報が含まれる。FECデータ部はFEC制御パケットにiSCSIネームが含まれる場合、データパケット、冗長パケットに存在する。

#### 【 0 0 3 0 】

図7はFECの制御テーブルと送信及び受信管理テーブルの内容の一例を示す図である。

## 【 0 0 3 1 】

F E C 許可テーブル119は F E C アダプタ毎に備えられており、通信相手が F E C 通信の対象か否かを管理するための所謂台帳として機能する。そのためこれには通信相手毎の iSCSI ネームと冗長度を登録する。このテーブルに登録されている iSCSI ネームに対するログインが行われたら、F E C 送信管理テーブル又は F E C 受信管理テーブルに登録される。尚、この許可テーブル119に登録されていない相手に対しては通常の T C P / I P 通信を行うようにして良い。

## 【 0 0 3 2 】

F E C 送信管理テーブル117には、個々の送信先アドレス（デスティネーションアドレス）と、それに対応した冗長度と、エンコード用バッファ制御情報が登録されている。ここで、デスティネーションアドレスとしては例えば I P アドレスが使用される。送信管理テーブル117にデスティネーションアドレスが登録されていれば、送信データである iSCSI データは F E C エンコード部113で F E C 処理されてネットワークへ送信される。即ち F E C 通信モードで送信される。一方デスティネーションアドレスが登録されていないければ、iSCSI データは F E C エンコード部113で F E C 処理されずに、T C P / I P の通信モードでネットワークへ送信される。

F E C 受信管理テーブル118には、ソースアドレスに対応してパケット到達率、F E C デコード用バッファ制御情報が登録される。ソースアドレスとしては受信先の I P アドレスが使用される。このパケット到達率からパケット喪失率が求められ、最終的には送信先への冗長度の変更のために反映される。

尚、これら F E C 送信管理テーブル及び受信管理テーブルには、更に iSCSI ネームの欄が追加されても良い。その欄には iSCSI ネームが登録され、F E C 許可テーブルからこの iSCSI ネームを参照することにより、これらのエンコード、デコード管理テーブルからの登録情報の追加、削除が行われる。

## 【 0 0 3 3 】

次に図 8 乃至図 1 0 を参照して、F E C 変換アダプタにおけるデータを送信するときの F E C 冗長化処理の動作について説明する。

## 【 0 0 3 4 】

## (A) エンコード処理

まず、図8に示すフローチャートを参照してエンコード処理について説明する。ストレージ装置で生成されたデータは、iSCSIパケット群として高速IPインターフェースを介してチャネルアダプタ107からFEC変換アダプタ110に送られる。FEC変換アダプタ110では物理層1101を介してこのパケット群が受信され(801)、バッファ1102に一時的に格納される。そして、FECエンコード部113では、パケットがiSCSI Login request PDUかまたはiSCSI Login response PDUかをチェックする(802)。その結果、否ならば、次にFEC送信管理テーブル117を参照し、デスティネーションアドレスがこの管理テーブル117登録されているか否かがチェックされる(803)。もし、登録されていれば、FECエンコード用バッファ1103が割り当てられる(804)。一方、制御テーブル117に登録されていなければ、パケットはFEC処理されずそのままの形式(即ちTCP/IP形式)のパケット群データとしてバッファ1103を介して物理層1104に送られ、ネットワークに送出される(816)。また、上記ステップ804で、FECエンコードバッファの領域が割り当てられていなければ、新しいFECシーケンスの開始準備を行う(817)。

## 【0035】

上記したFEC用エンコードバッファが割り当てられる場合、パケットを格納するアドレスが計算され(805)、パケットはそのアドレスへそのまま格納される(806)。即ちIPパケットは丸ごとカプセル化され、IPヘッダ、TCPヘッダ、TCPデータはそのままバッファに格納される。次に、パケットにはIPヘッダ、UDPヘッダ、FECヘッダが付加され、カプセル化して、物理層1104を介してネットワークへ送出される(807)。FECヘッダの各領域には対応する値を格納する。即ちFECヘッダ種にはFECデータを、データ長にはカプセル化したパケットの長さを格納し、その他はFEC送信管理テーブルの内容に従う。

その後、FEC送信管理テーブル117の内容が更新される(808)。このテーブル117の更新はこのFECシーケンスについて格納されるパケット数をプラス“1”する。送信するパケットの総数を管理するためにある。

## 【 0 0 3 6 】

次に、F E Cシーケンス（冗長度  $n$ ）に対して  $n$  個目のパケットか否かが判断され（809）、 $n$  個目でなければ終了する。一方、 $n$  個目であれば、最後に冗長データを  $n+1$  個目のパケットして送信するため、冗長度  $n$  に従ってこの F E C エンコード用バッファのデータから冗長データを作成する（810）。そして、冗長データに I P ヘッダ、U D P ヘッダ、F E C ヘッダを付加し、カプセル化して物理層1104を介してネットワークに送出される（811）。そしてこの一連の F E C 処理シーケンスの終了処理をする（812）。終了処理とは、F E C 再送機能が有る場合には、この F E C シーケンスは F E C \_ A C K 待ちとするものであり、F E C 再送機能が無い場合には、F E C エンコードバッファを開放するなどの処理である。

さて、上記ステップ 8 0 2 でのチェックの結果、Yes ならば、そのパケット群が i S C S I Initial Logon request P D U か否かのチェックが行われる（813）。その結果 Yes であれば、ターゲットネームが F E C 許可テーブルに登録されているか否かがチェックされる（814）。このテーブル登録のチェックの結果、登録されていれば、パケットのデスティネーションアドレスについて F E C 送信管理テーブルに追加して次の処理へと移る（815）。

一方、上記ステップ 8 1 3 のチェックの結果、否の場合にはそのパケットを F E C 処理せずそのまま物理層1104へ送り（816）、処理を終了する。

## 【 0 0 3 7 】

## （B）送信データが不足している場合の処理動作

次に、図 9 に示すフローチャートを参照して、送信データが不足している場合の処理動作について説明する。この処理はストレージ装置のチャネルアダプタから正規の数のパケットを受け取っているか否かを確認するための処理である。

ある送信シーケンスにかかるパケット群の送信を開始すると、タイマ割り込みがかけられ（820）、タイマのカウントが開始される。そしてタイマが、送信処理中の F E C シーケンスの処理開始から一定時間以上経過したか否かがチェックされる（821）。経過していなければ処理は終了する。

一方、一定時間を経過していれば、この F E C シーケンスについて既に物理層11

01から受信しているm個のパケットに対して、このFECシーケンスの冗長度をmとして冗長データを作成する(822)。尚、この冗長度mは冗長パケットのFECヘッダに格納される。

#### 【0038】

##### (C) FECパケットの再送動作

次に、図10のフローチャートを参照してFECパケットの再送動作について説明する。

この処理の基本的シーケンスは、受信側ではFECパケットを受信するとFEC\_\_ACKを送信側に返送する。送信側ではこのFEC\_\_ACKを受領することにより、送信したパケットが正しく相手に送られたと認識する。

さて、送信側ではFECパケットを送信すると、タイマ割り込み処理が行われる(830)。そしてあるFEC\_\_ACK待ち中のFECシーケンスに関してタイマがタイムアウトをしたか否かがチェックされ(831)る。このチェックの結果タイムアウトしていなければ終了する。しかし、チェック831の結果、タイムアウトした場合には、このFECシーケンスについての全パケットを再送信するための処理が行われる。同様に、これら再送のパケットに関してFEC\_\_ACKを待つことになるので、タイマも再設定される(832)。

#### 【0039】

次に、図11及び図12を参照してFEC変換アダプタ110でデータを受信するときのFECデコード処理の動作について説明する。

図11及び図12に示すフローチャートにおいて、物理層1114からパケットを受信すると(901)、パケットがiSCSI Login request PDU 又はiSCSI Login response PDUであるか否かがチェックされる(902)。Yesであれば、そのパケットがiSCSI initial Login request PDUか否かがチェックされ(919)、その結果該当すれば、source ネーム(即ちiSCSI ネーム)がFEC許可テーブルに登録されているか否かがチェックされる(920)。許可テーブルに登録されていれば、そのパケットのデスティネーションアドレスについてFEC受信管理テーブルに登録する(921)。

上記チェック902の結果、否であればデスティネーションポート番号がFEC

通信用か否かがチェックされる (903)。このチェック 9 0 3 で、否の場合にはそのパケットを F E C 処理しないでそのまま物理層 1111 へ送る (922)。前記デスティネーションポート番号のチェック 9 0 3 の結果、F E C 通信用であれば、F E C 許可テーブル情報の変更パケットか否かがチェックされる (904)。その結果、Yes であれば、F E C 許可テーブルの内容を更新して (923)、T C P 処理を行う (924)。即ち A C K を返信して、処理を終わる。

#### 【 0 0 4 0 】

一方、T C P で無ければ、U D P パケットか否かがチェックされ (905)、U D P パケットであれば、ソースアドレスが F E C 受信管理テーブルに在るか否かがチェックされる (906)。ソースアドレスがこの受信管理テーブルに在れば、パケットが破損しているか否かをチェックし (907)、そしてタイムアウトになった F E C シーケンスのパケットか否かがチェックされる (908)。これらのチェックの結果、Yes であれば、F E C 受信管理テーブルについてソースアドレスからパケット到着率を更新処理し (925)、そのパケットを廃棄して終わる (926)。

一方、前記チェック 9 0 8 によりタイムアウトになった F E C シーケンスのパケットでなければ、F E C 受信管理テーブルについてソースアドレスからパケットの到達率を更新処理する (909)。その後、デコード済みの F E C シーケンス I D を持つか否かがチェックされ (910)、Y e s であれば、そのパケットを破棄して終わる (926)。デコード済みの F E C シーケンス I D で無ければ、続いて、その F E C シーケンスについて F E C デコード用バッファの領域が割り当てられているか否かがチェックされる (911)。その結果否ならば、F E C シーケンスのデコード処理を開始する (927)。これは例えば F E C デコード用バッファの未使用領域を割り当てて、タイマセットすることにより行われる。

#### 【 0 0 4 1 】

ステップ 9 1 1 のチェックで、該当すれば、そのパケットを格納するアドレスを計算して (912)、そのパケットを該当するアドレスへ格納する (913)。次に、そのパケットが F E C 制御パケットか否かをチェックする (914)。チェックの結果、否であれば、そのパケットがデータパケットか否かをチェックし (

928)、データパケットであれば、そのデータをカプセル化したパケットとして物理層1111に送る(929)。

次にFECシーケンス(この場合冗長度は $n$ )についてそのパケットが $n$ 個目に届いたか否かをチェックする(915)。その結果該当すれば、続いてFECシーケンスについて未到着のパケットは冗長パケットか否かがチェックされる(916)。冗長パケットであれば、FEC\_ACKを返信して、シーケンスIDを格納する(917)。そしてこの受信FECシーケンスの処理を終了する(918)。即ち、FECデコード用バッファの割り当てを開放して、シーケンスIDを処理済みシーケンスとしてFEC受信管理テーブルに一時的に格納する動作をする。

#### 【0042】

上記ステップ916における冗長パケットか否かのチェックで、該当しなければ、冗長度 $n$ (ステップ822により冗長パケットのFECヘッダが冗長度 $m$ を示すならば、 $m$ )に従って未到着のパケットを再構成する(930)。再構成されたパケットがFEC制御パケットか否かがチェックされ(931)、該当しなければ、再構成したパケットがFECデータパケットか否かをチェックする(932)。そしてFECデータパケットであれば、再構成したデータをカプセル化したパケットとして物理層1111へ送る(934)。尚、FEC制御パケットの処理は後述する。

#### 【0043】

次に図13乃至図15を参照して、データの受信時のエラーが発生した場合の処理、及びパケット喪失に係わる処理動作等について説明する。

#### 【0044】

##### (A) 受信データ不足時の処理

この処理はFEC冗長化されたFECシーケンスが一定時間の内に符号化できるほど十分に到達しなかった場合の処理である。

図13に示すフローチャートにおいて、あるシーケンスのパケットの受信を開始すると、タイマ割り込みをする(1501)。そしてあるシーケンス(冗長度 $n$ )について、処理を開始してから受信パケットが $n$ 個未満のままで一定時間が経過したか否かがチェックされる(1502)。 $n$ 未満であれば、FECシーケンスの未到



達パケットの数 $x$ を算出する(1503)。これは例えば、(n-到達パケット数)の計算をすることにより算出可能である。

次にFECシーケンスを格納したデコード用バッファ内をクリアして(1504)、パケット到達率の更新を行って、終わる(1505)。即ちこのFECシーケンスのうち、未到達のパケットは喪失したものとして扱い、パケット到達率を計算する。その計算結果が更新されたパケット到達率となる。

#### 【0045】

##### (B) パケットの到達率の報告処理

算出されたパケット到達率は受信側では受信管理テーブルに登録されると共に、送信側にも送られ、そこでパケット喪失率から冗長度が求められ、冗長度のデータとして送信管理テーブルに登録される。

図14に示すフローチャートにおいて、受信側ではタイマ割り込みを行い(1601)、各FECソースアドレスに対して、パケット到達率を報告するための管理フレームを送信する(1602)。パケット到達率報告のパケットはFECエンコード用バッファに格納される。パケット到達率(Num\_Alive PKT)は、破損、タイムアウトしないで届いたFECデータパケット数と、送信されたFECデータパケット数の比を取るにより算出される(1603)。

#### 【0046】

##### (C) FEC制御パケットの処理

図15に示すフローチャートにおいて、送信側では、届いたパケットがパケット喪失率の報告パケットか否かをチェックする(1701)。喪失率の報告パケットであれば、報告されたパケット喪失率と、予め決められた安全係数から、そのデスティネーション宛ての冗長度を算出して決められる(1702)。求められた冗長度は送信管理テーブルの該当する個所を更新する形で登録される。また、上記ステップ1701で、パケット喪失率報告用のパケットで無ければ、管理テーブル変更のパケットか否かがチェックされる(1704)。このチェック1704で、該当すれば、送信管理テーブルの内容が更新される(1705)。一方、管理テーブル変更パケットでもなければ、FEC\_ACKか否かがチェックされる(1706)。その結果該当すれば、送信FECシーケンスの処理を終了する(1707)。この処理

はFECエンコード用バッファの割り当て領域を開放し、FEC\_\_ACK待ちを解除し、FEC送信管理テーブルから削除する処理である。

## 【0047】

以上、本発明の実施例を説明したが、上記実施例に限定されずに、種々変形して実施し得る。

## 【0048】

図16はストレージ装置のチャネルアダプタの構成を変形した例である。このチャネルアダプタ107に、図2を参照して述べたFEC変換アダプタの持つ機能を併合したものである。即ち、図2に示した構成のうち、チャネルアダプタ側にある構成、バッファ1102、1112、及び物理層1101、1111が省略され、チャネルアダプタ107はSCSIコマンド制御部1071、プロトコル制御部1072を有しているが、送信バッファ1073、受信バッファ1074が兼用される。残余の構成は図2と同様である。

## 【0049】

また、上記実施例ではパケットの喪失率に従って、送信先毎に冗長度を変更してデータを送信するようにしたが、パケットの喪失率に従ってパケット送信の間隔を変えることも可能である。広義にはパケット喪失率に従ってデータ送信の状態を変えるようにできる。

## 【0050】

## 【発明の効果】

本発明によれば、FEC方式を採用したiSCSIプロトコル準拠のストレージ装置を提供することができる。また、受信側で把握したパケットの喪失状況を送信側にフィードバックすることにより、送信側ではデータの送信先毎にFEC処理の状態やデータ送信の冗長度を変えてデータ送信することができ、一方受信側ではそれに応じて復元することができる。これによりパケットが喪失した場合でもデータの復元が可能はストレージシステムを得ることができる。

## 【図面の簡単な説明】

## 【図1】

本発明の一実施例によるストレージシステムを示すブロック図。

【図 2】

本発明の一実施例に係る F E C 変換アダプタを有するストレージ装置の一例を示す図。

【図 3】

本発明の他の実施例によるストレージシステムを示すブロック図。

【図 4】

本発明の一実施例によるデータ転送の概念を示す図。

【図 5】

通信に用いられるパケットのフォーマットの一例を示す図。

【図 6】

F E C ヘッダのフォーマットを示す図。

【図 7】

F E C の制御テーブルと送信及び受信管理テーブルの内容の一例を示す図。

【図 8】

データを送信するときの F E C エンコード処理の動作を説明するためのフローチャート。

【図 9】

F E C エンコード処理における送信データ不足時の動作を説明するためのフローチャート。

【図 1 0】

F E C パケットの再送処理動作を説明するためのフローチャート。

【図 1 1】

データ受信時における F E C デコード処理の動作を説明するためのフローチャート。

【図 1 2】

データ受信時における F E C デコード処理の動作を説明するためのフローチャート。

【図 1 3】

データの受信時にエラーが発生した場合の処理動作を説明するためのフローチャート。

ート。

【図 1 4】

データの受信時におけるパケット到達率の報告処理動作を説明するためのフローチャート。

【図 1 5】

パケット喪失の報告に対する処理動作について説明するためのフローチャート。

【図 1 6】

本発明の他の実施例によるストレージ装置の構成を示す図。

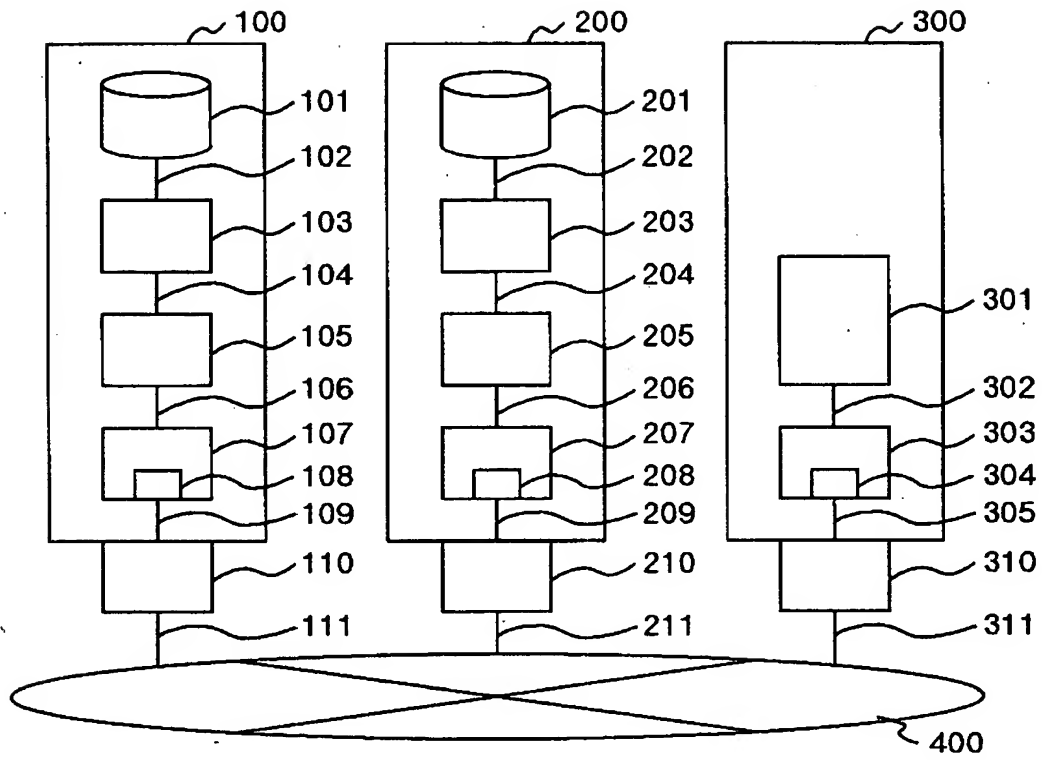
【符号の説明】

100,200:ストレージ装置	300:ホストコンピュータ
400:ネットワーク	500:管理サーバ
101,201:ディスクドライブ	103,203:ディスクアダプタ
105,205:キャッシュメモリ	107,207:チャネルアダプタ
110,210,310:F E C変換アダプタ	

【書類名】 図面

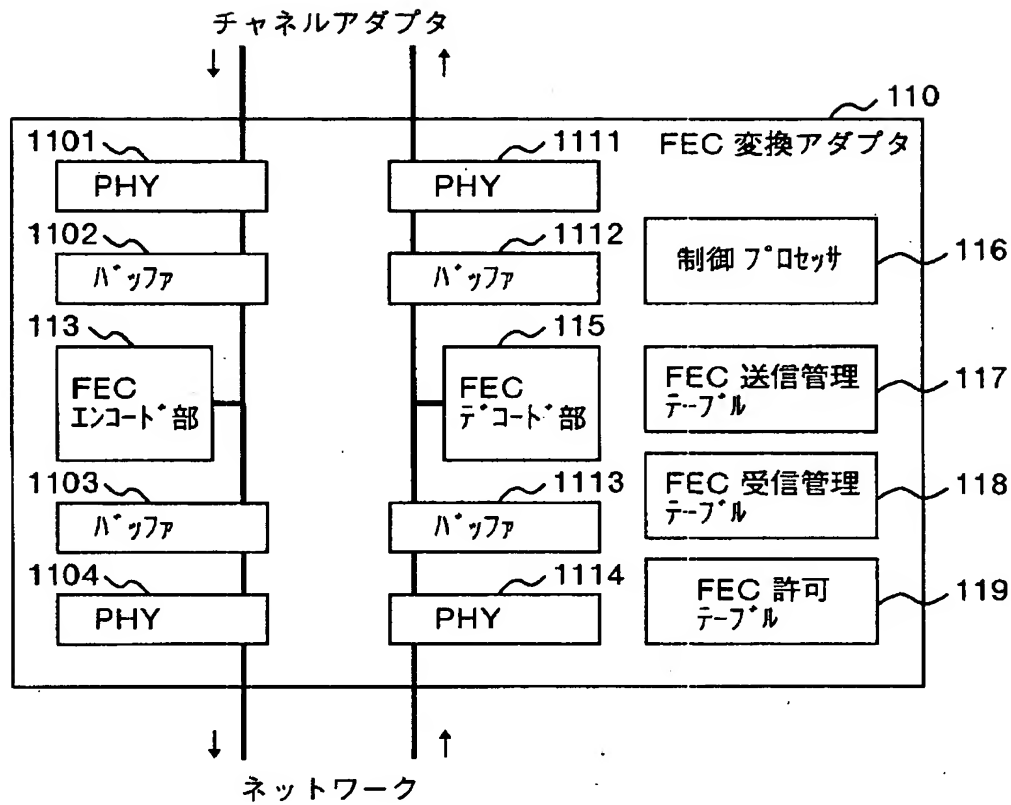
【図 1】

図 1



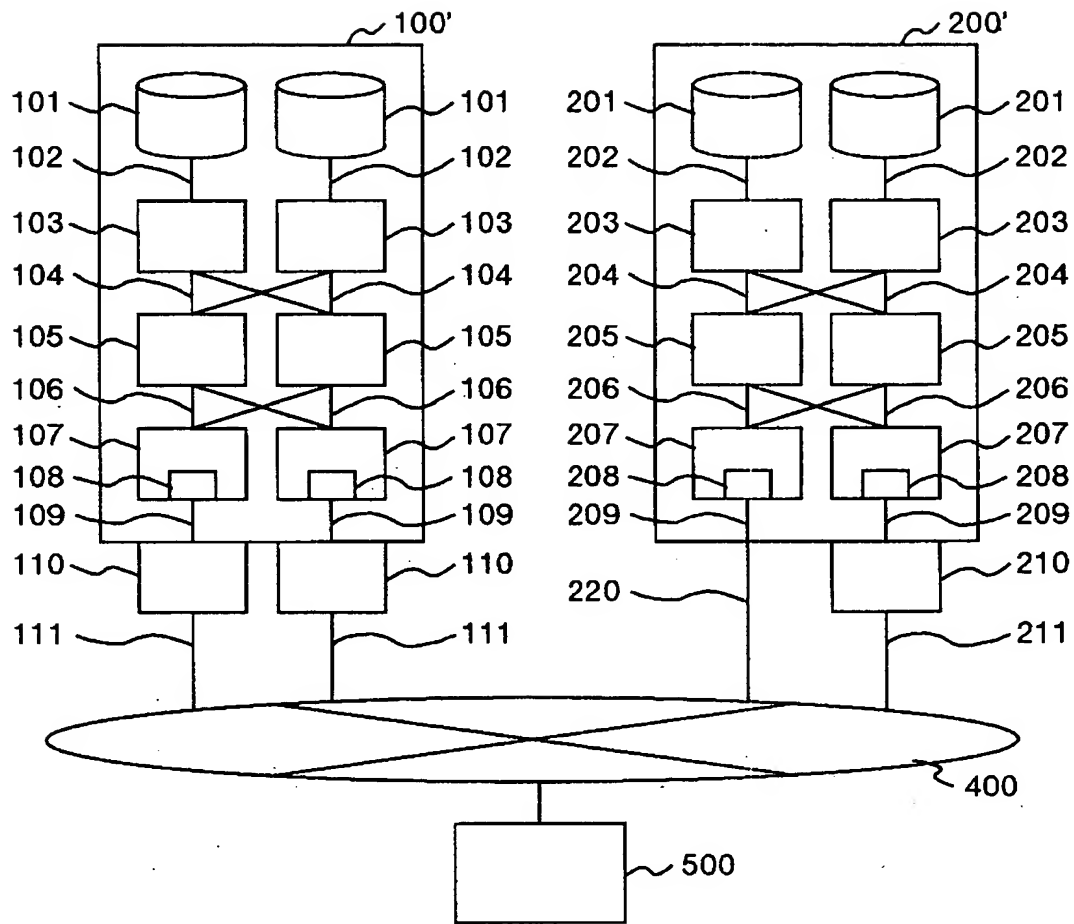
【図 2】

図 2



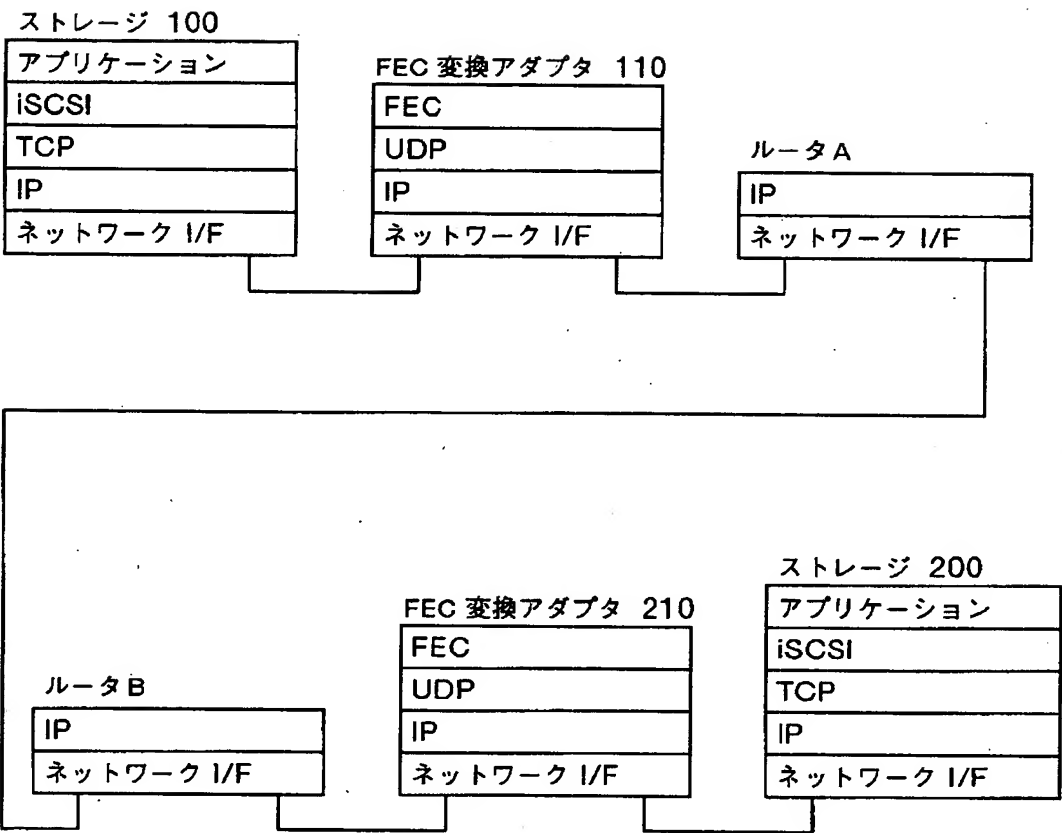
【図 3】

図 3



【図 4】

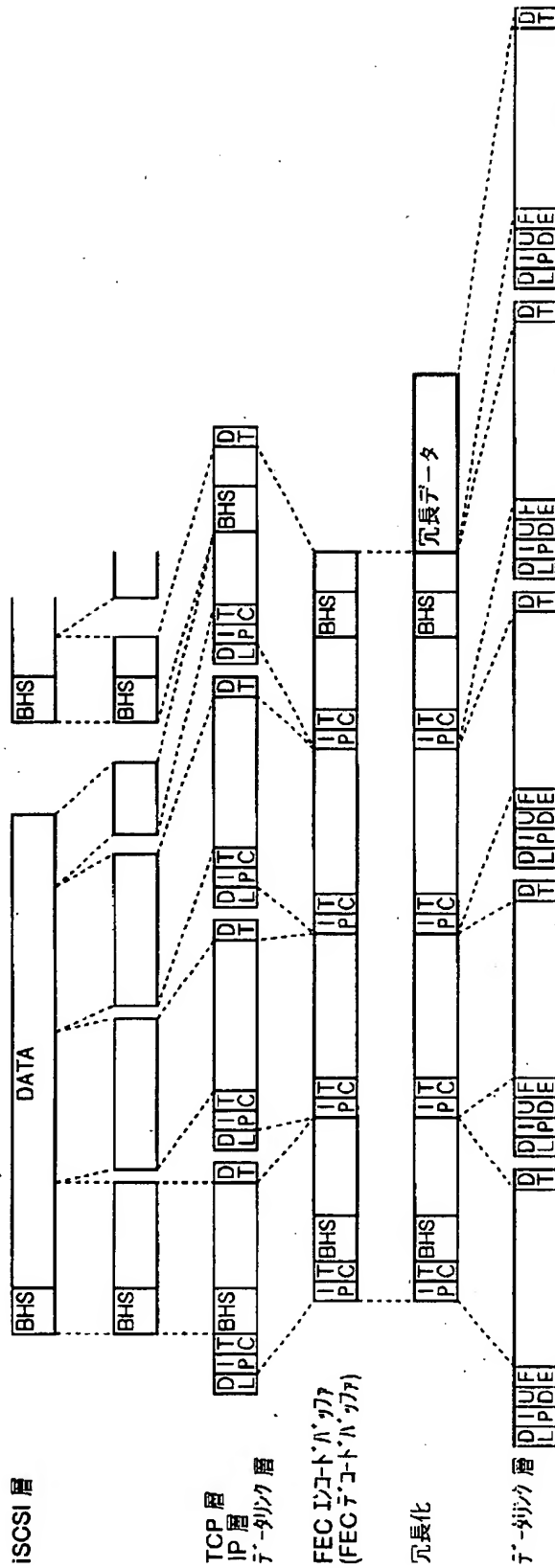
図 4





【図 5】

図 5

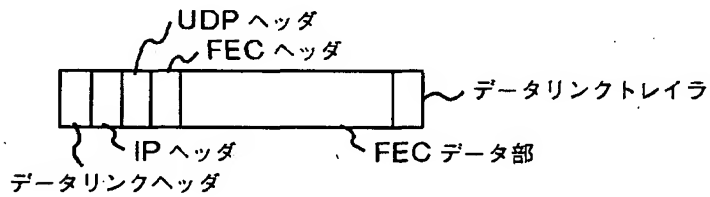


DL : データリンク層  
 IP : IP 層  
 TC : TCP 層  
 DT : データ転送層  
 UD : UDP 層  
 FE : FEC 層

【図 6】

図 6

ビット																																			
0										1										2										3					
0 1 2 3 4 5 6 7 8 9										0 1 2 3 4 5 6 7 8 9										0 1 2 3 4 5 6 7 8 9										0 1					
ワード 0	エンコード情報										制御情報種															FEC 制御情報 ID									
ワード 1	FEC 制御情報 (制御情報種による)																																		
ワード 2	FEC 制御情報 (制御情報種による)																																		
ワード 3	FEC 制御情報 (制御情報種による)																																		



【図 7】

図 7

FEC 許可テーブル 119

通信相手の iSCSI Name	初期冗長度 (n)

FEC 送信管理テーブル 117

デスティネーションアドレス	冗長度 (n)	FEC エンコード用バッファ制御情報

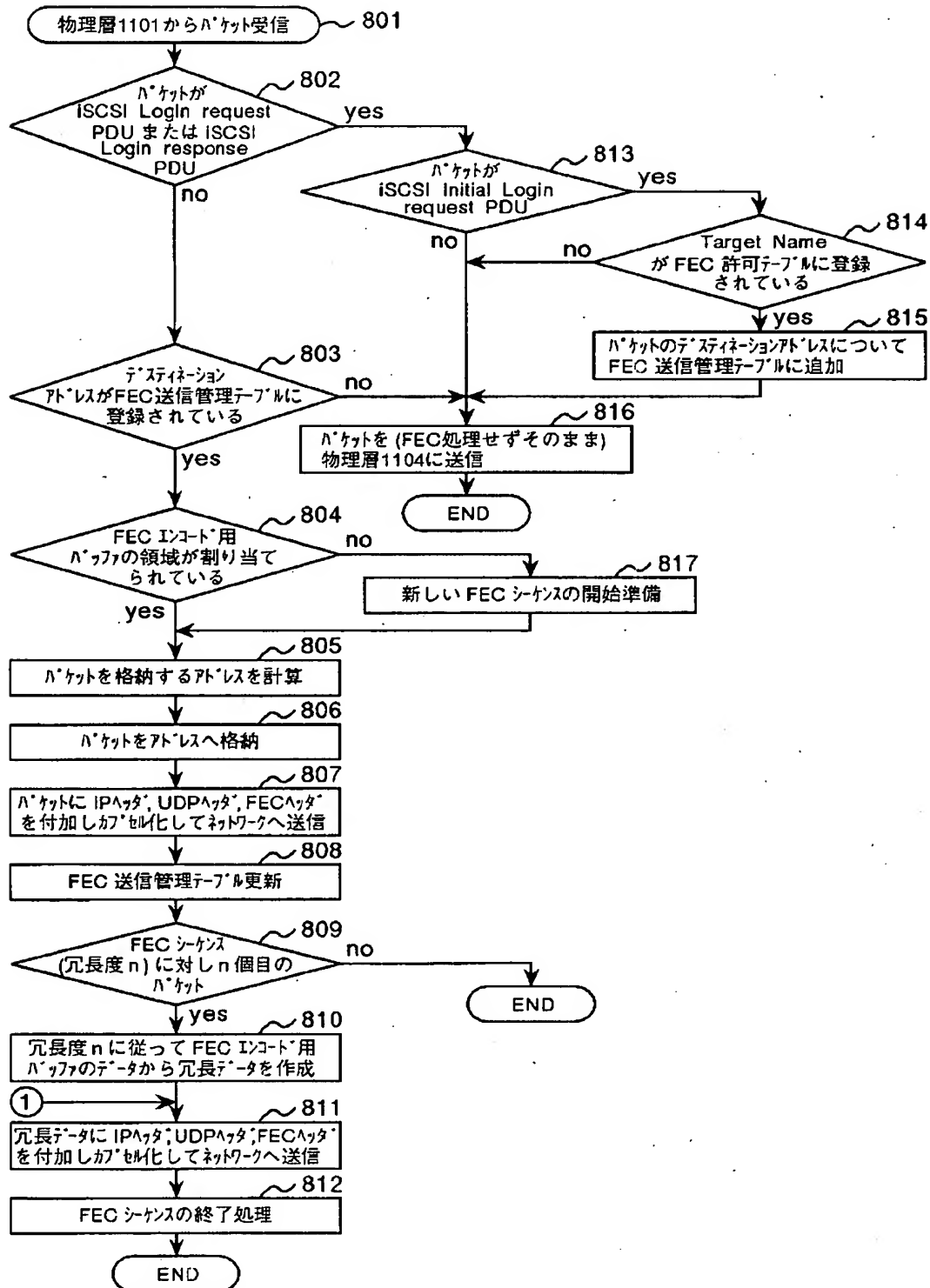
FEC 受信管理テーブル 118

ソースアドレス	パケット到達率	FEC デコード用バッファ制御情報
	/	
	/	
	/	
	/	

【図 8】

図 8

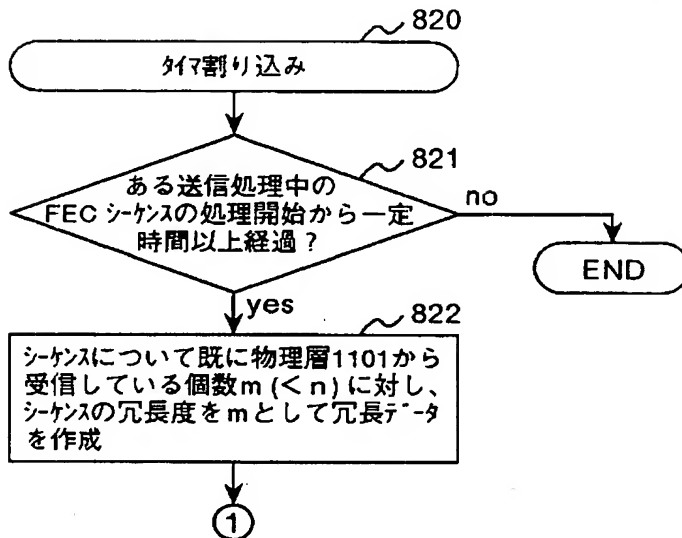
エンコード処理



【図 9】

図 9

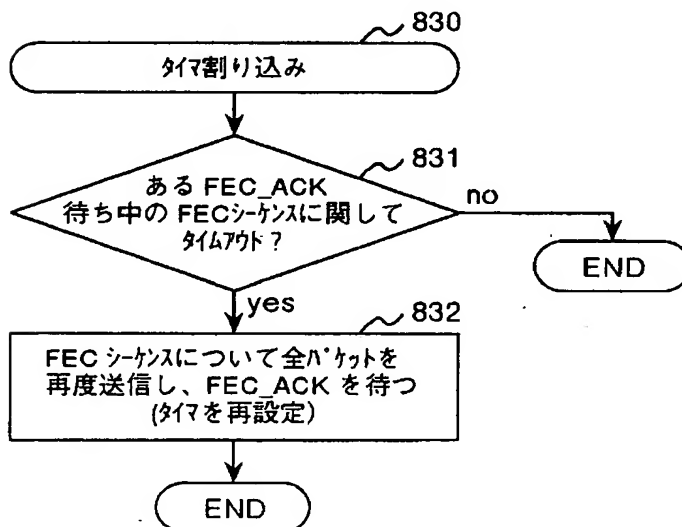
送信データ不足時の処理



【図 1 0】

図 1 0

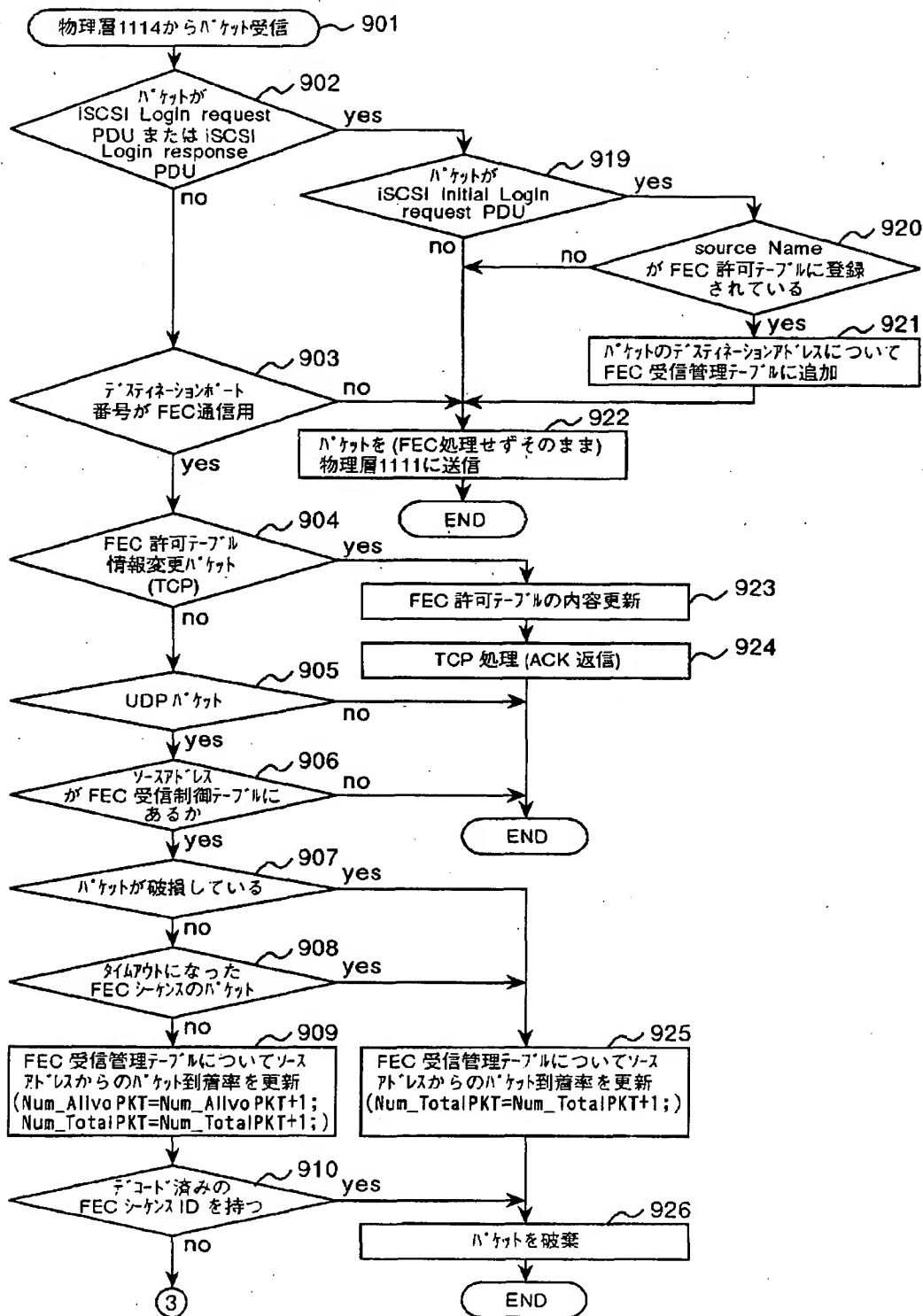
FEC\_ACK 待ちタイムアウトによる再送処理



【図 11】

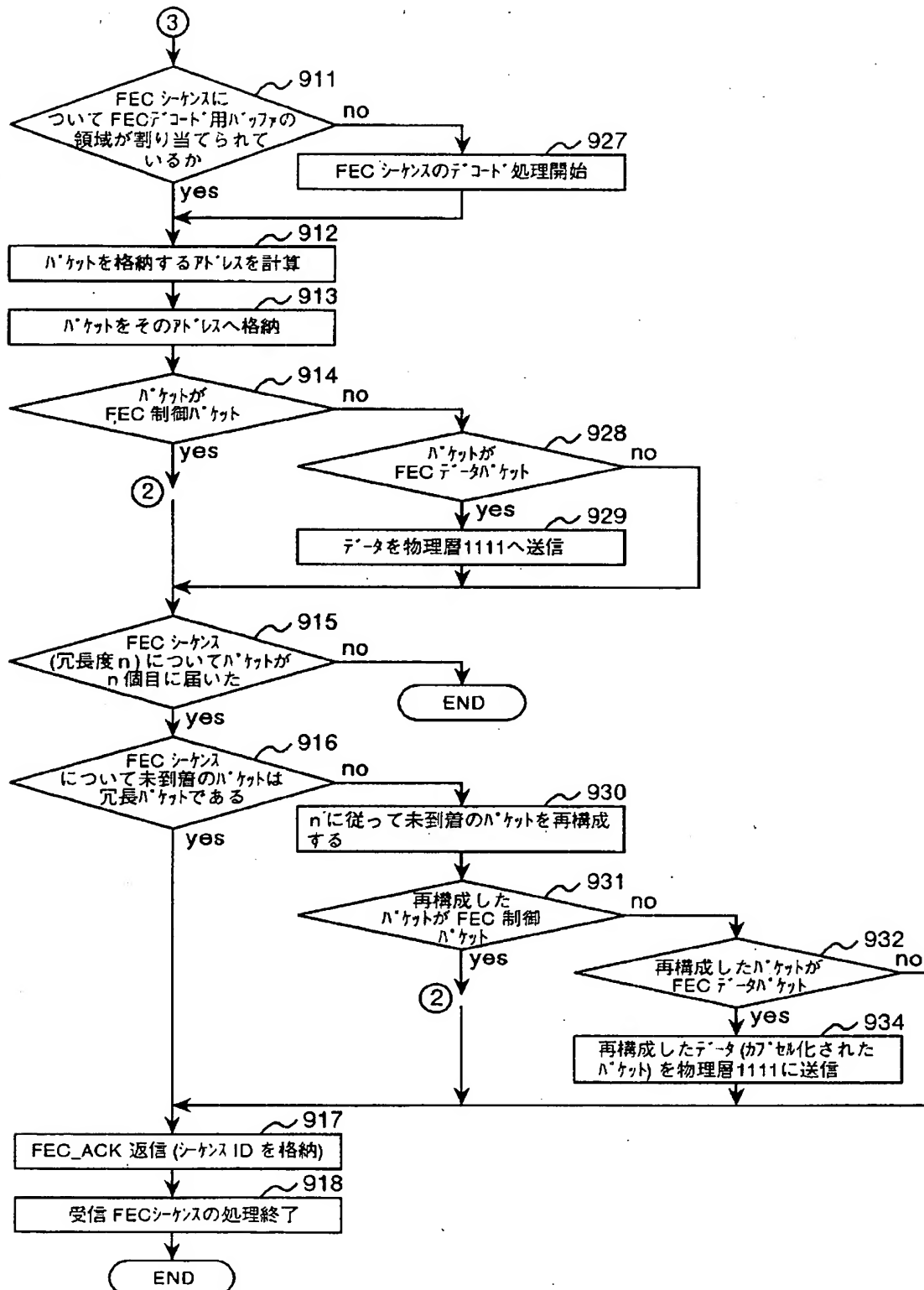
図 11

FEC デコード処理



【図 12】

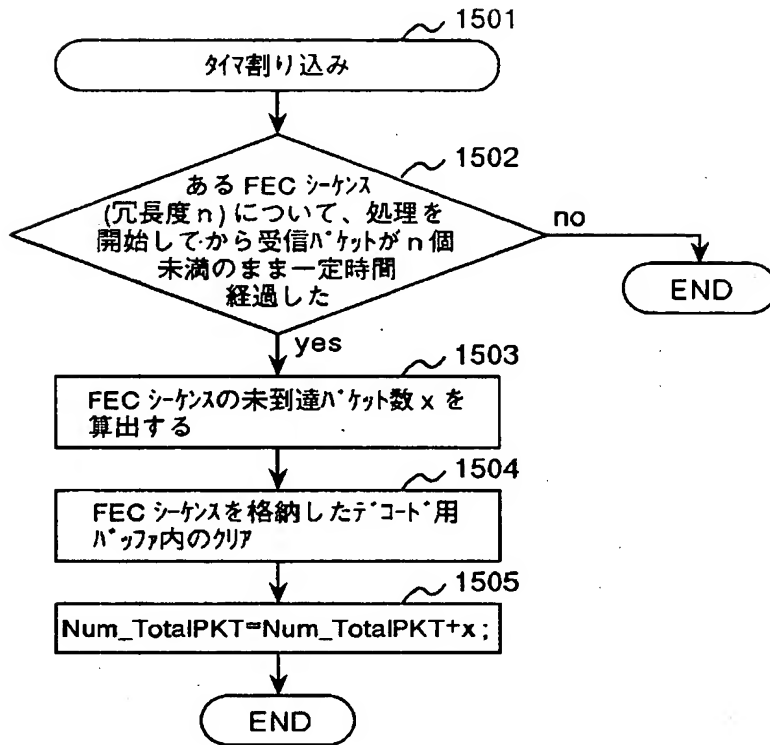
図 12



【図 1 3】

図 1 3

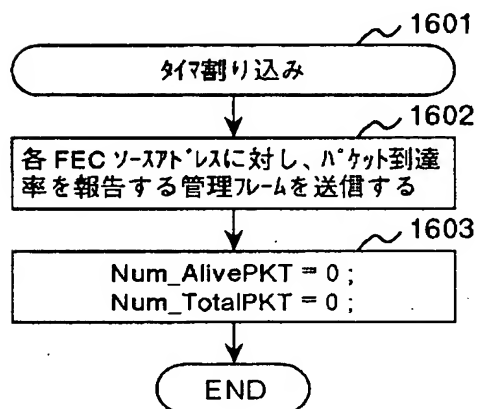
受信データ不足時の処理



【図 1 4】

図 1 4

パケット到達率の報告

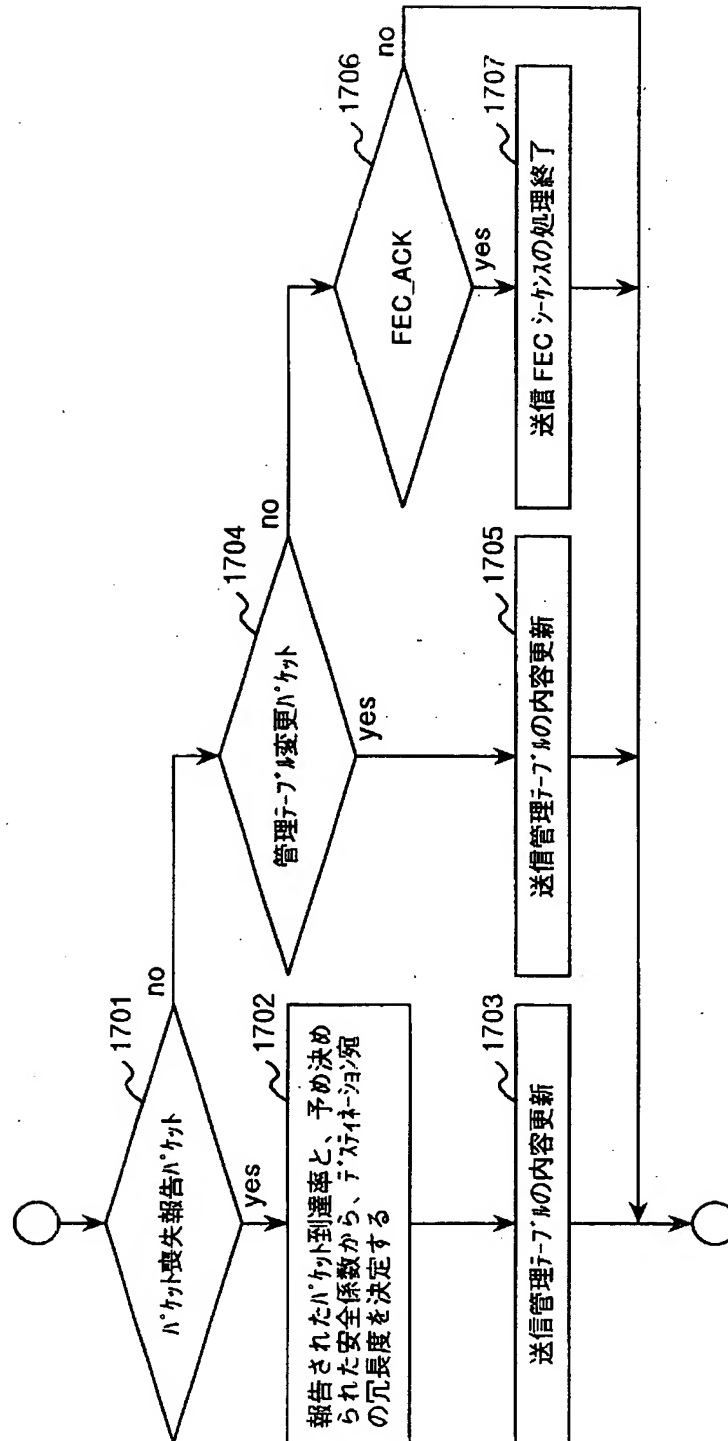




【図15】

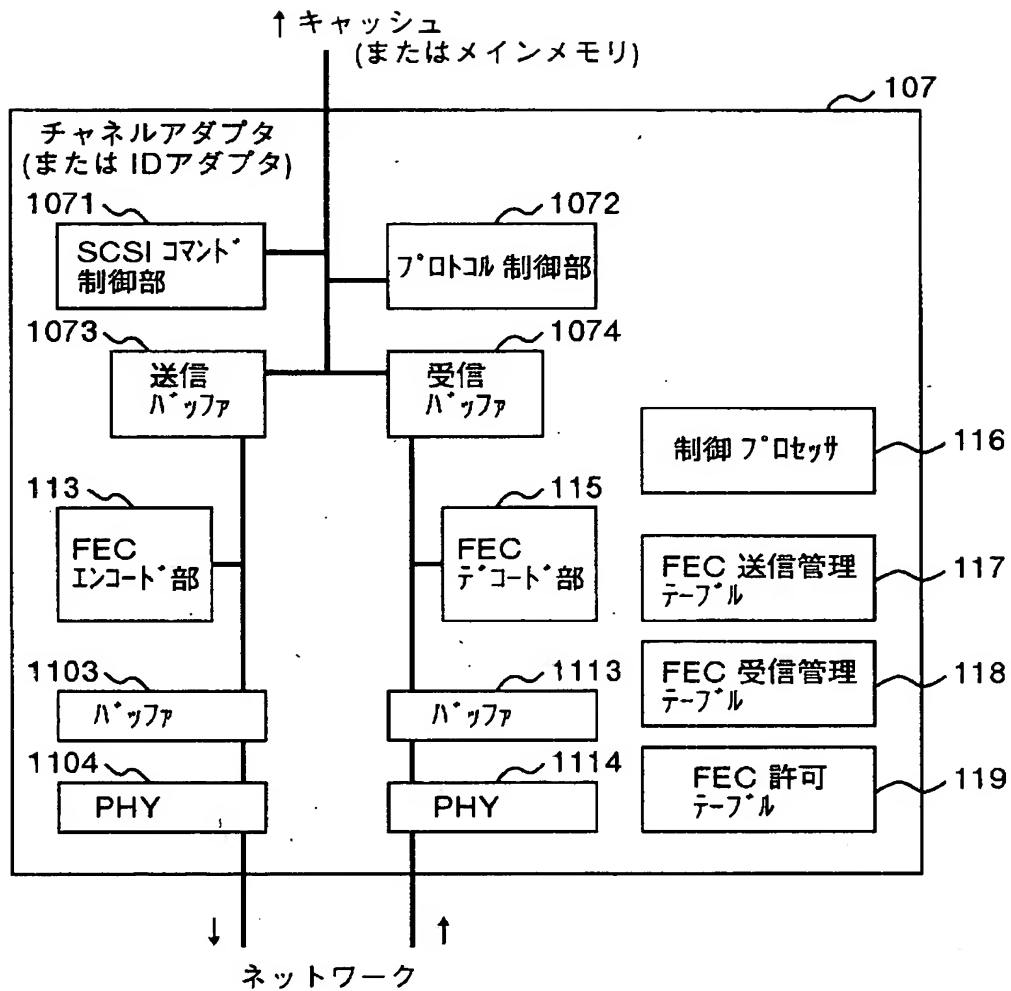
図 15

FEC 制御パケット処理



【図 16】

図 16



【書類名】 要約書

【要約】

【課題】

FEC方式を採用したiSCSIプロトコル準拠のストレージ装置間のデータ送受信において、パケットが喪失した場合でも送信データの復元が可能なストレージシステムを実現する。

【解決手段】

iSCSIプロトコル準拠のストレージ装置がネットワークを介して接続されたストレージシステムにおいて、ストレージ装置にFECの機能を持たせ、iSCSI層間でデータの伝送を行う場合、データの送信先毎にFEC処理の状態やデータ送信の冗長度を変えてデータ送信し、復元する。

例えば、ネットワーク側のポートとストレージ装置側のポートとの間に介在し、パケットデータの送受信を行うシステムにおいて、送信先に対応してFEC冗長度を登録して管理する送信管理テーブルと、受信先に対応してFEC冗長度を登録して管理する受信管理テーブルと、ストレージ装置で生成されるパケット化されたiSCSI層のデータを、該送信管理テーブルを参照して送信先に応じた冗長度を持たせてFECエンコード処理するエンコード部と、該受信管理テーブルを参照して、ネットワークから受信されるパケットデータをFECデコード処理し、iSCSI層のデータに復号化するデコード部とを有する。

【選択図】 図1

出 願 人 履 歴 情 報

識別番号 [0 0 0 0 0 5 1 0 8]

1. 変更年月日	1 9 9 0 年 8 月 3 1 日
[変更理由]	新規登録
住 所	東京都千代田区神田駿河台 4 丁目 6 番地
氏 名	株式会社日立製作所